

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 07-104791

(43)Date of publication of application : 21.04.1995

(51)Int.Cl.

G10L 7/02

(21)Application number : 05-247828

(71)Applicant : ATR ONSEI HONYAKU TSUSHIN
KENKYUSHO:KK

(22)Date of filing : 04.10.1993

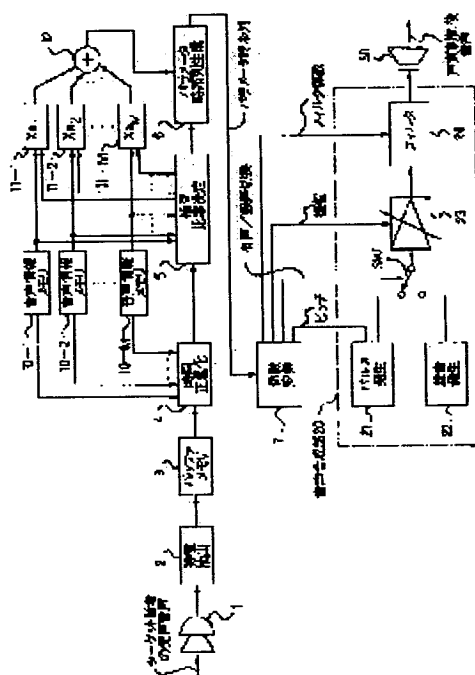
(72)Inventor : IWAHASHI NAOTO
KOSAKA YOSHINORI

(54) VOICE QUALITY CONTROL TYPE VOICE SYNTHESIZING DEVICE

(57)Abstract:

PURPOSE: To obtain a device which can synthesize voice after voice quality control to voice inputted with a higher degree of freedom through a simple constitution by performing an interpolating process for data on voice spectra of plural speakers by using a specific interpolation rate and outputting data on a voice spectrum having specific voice quality, and synthesizing the voice on the basis of the data.

CONSTITUTION: Voice information memories 10-1-10-M store the data of the voice spectra of plural M speakers in advance respectively, the voice of a target speaker is inputted, and the voice spectra of the M speakers are interpolated and mixed so as to approximate the input voice, thereby vocalizing a voice close to the voice quality of the voice of the target speaker. And, a voice having specific different quality is vocalized by varying the interpolation ratio. Here, a cepstrum coefficient, etc., having relatively excellent interpolation characteristics is used as a spectrum parameter for interpolation, and the interpolation ratio of the parameter is varied to synthesize and output a voice having different quality.



BEST AVAILABLE COPY

LEGAL STATUS

[Date of request for examination]

04.10.1993

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]

[Patent number]

2951514

[Date of registration]

09.07.1999

[Number of appeal against examiner's decision]

THIS PAGE BLANK (USPTO)

of rejection]

[Date of requesting appeal against examiner's
decision of rejection]

[Date of extinction of right]

09.07.2003

THIS PAGE BLANK (USPTO)

(19)日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11)特許出願公開番号

特開平7-104791

(43)公開日 平成7年(1995)4月21日

(51)Int.Cl.⁶

G 1 0 L 7/02

識別記号

庁内整理番号

D

F I

技術表示箇所

審査請求 有 請求項の数 4 O L (全 6 頁)

(21)出願番号 特願平5-247828

(22)出願日 平成5年(1993)10月4日

(71)出願人 593118597

株式会社エイ・ティ・アール音声翻訳通信
研究所

京都府相楽郡精華町大字乾谷小字三平谷5
番地

(72)発明者 岩橋 直人

京都府相楽郡精華町大字乾谷小字三平谷5
番地 株式会社エイ・ティ・アール音声翻
訳通信研究所内

(72)発明者 匂坂 芳典

京都府相楽郡精華町大字乾谷小字三平谷5
番地 株式会社エイ・ティ・アール音声翻
訳通信研究所内

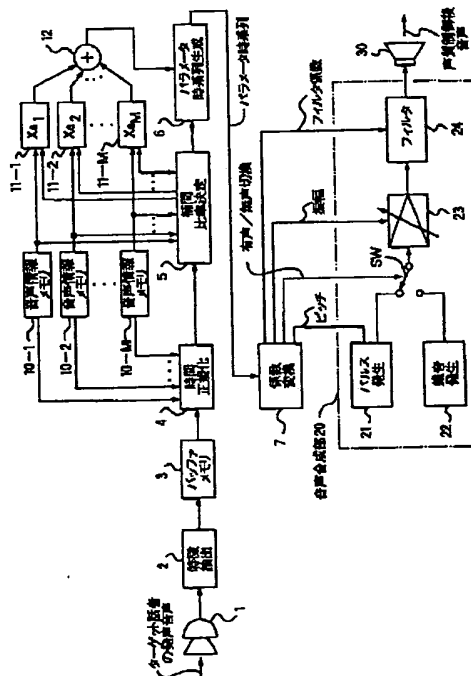
(74)代理人 弁理士 青山 葆 (外2名)

(54)【発明の名称】 声質制御型音声合成装置

(57)【要約】

【目的】 従来例に比較して非常に簡単な構成を有し、しかもより高い自由度を有して入力された音声に対して声質制御した後の音声を合成することができる音声合成装置を提供する。

【構成】 複数の話者の音声スペクトルのデータを予め記憶する記憶装置と、記憶装置から複数の話者の音声スペクトルのデータを読み出し、所定の補間比率を用いて上記複数の話者の音声スペクトルのデータに対して内挿処理を実行して所定の声質を有する音声スペクトルのデータを出力する処理回路と、処理回路から出力される音声スペクトルのデータに基づいて音声を合成して出力する音声合成回路とを備える。さらに、ターゲット話者の音声に基づいて音声スペクトルのデータを抽出し、抽出された音声スペクトルのデータが処理回路から出力される音声スペクトルのデータに近似するように補間比率を演算して設定する演算回路とを備える。



【特許請求の範囲】

【請求項1】 複数の話者の音声スペクトルのデータを予め記憶する記憶手段と、

上記記憶手段から複数の話者の音声スペクトルのデータを読み出し、所定の補間比率を用いて上記複数の話者の音声スペクトルのデータに対して内挿処理を実行して所定の声質を有する音声スペクトルのデータを出力する処理手段と、

上記処理手段から出力される音声スペクトルのデータに基づいて音声を合成して出力する音声合成手段とを備えたことを特徴とする声質制御型音声合成装置。

【請求項2】 声質制御型音声合成装置はさらに、ターゲット話者の音声に基づいて音声スペクトルのデータを抽出する特徴抽出手段と、

上記特徴抽出手段によって抽出された音声スペクトルのデータが上記処理手段から出力される音声スペクトルのデータに近似するように上記補間比率を演算して上記処理手段に設定する演算手段とを備えたことを特徴とする請求項1記載の声質制御型音声合成装置。

【請求項3】 声質制御型音声合成装置はさらに、上記補間比率を入力して上記処理手段に設定する入力手段を備えたことを特徴とする請求項1記載の声質制御型音声合成装置。

【請求項4】 上記音声スペクトルのデータは、ケプストラム係数又は対数面積比であることを特徴とする請求項1、2又は3記載の声質制御型音声合成装置。

【発明の詳細な説明】

【0001】

【産業上の利用分野】 本発明は、声質を制御して音声合成することができる声質制御型音声合成装置に関する。

【0002】

【従来の技術及び発明が解決しようとする課題】 男性の声から女性の声に又は女性の声から男性の声に変換する方法（以下、第1の従来例という。）が、例えば、箱田和雄，“極制御による男女声変換法の検討”，日本音響学会講演論文集，2-6-13，pp213-214，昭和62年10月に開示されている。この第1の従来例においては、LPC分析で得られる極周波数を用いて、母音情報の代わりに極周波数値を用いて極の変換を行うとともに、変換に伴うスペクトルの傾きの変動を2次のフィルタを用いて吸収する方法を提案している。しかしながら、男性の声から女性の声に又は女性の声から男性の声への変換のみで、例えば男性から男性への変換で、異なる個性を有する音声を発生することはできない。

【0003】 また、アナウンサーなどの明瞭は音声の物理的特性の1つとして、ホルマント周波数のダイナミックな変動に着目し、特に声の明瞭性を改善するために、ホルマント周波数の時間変化を制御する方法（以下、第2の従来例という。）が、例えば、都木徹，桑原尚夫，“ホルマント変化の強調・抑圧による声質制御”，日本

音響学会講演論文集，1-4-12，pp145-146，昭和61年10月に開示されている。この第2の従来例は、次のステップを有する。

(a) 所定の標準化周波数でD/A変換した音声信号を、所定のフレーム幅、フレーム周期、及び男女で異なる分析次数で線形予測分析し、予測係数との残差を計算する。

(b) 各フレーム毎に予測係数からホルマント周波数を算出し、従来の粕谷の方法（“線形予測分析法で得られる極周波数からのホルマント周波数選択アルゴリズム”，電子通信学会論文誌，Vol. J66-A，No. 11，pp. 1144-1145，1983年11月参照。）により、母音部の第1乃至第3ホルマントの軌跡を求める。

(c) 得られたホルマント軌跡に対して、ある音声のホルマント周波数の時間変化を示す所定の式に適用し、各フレーム毎に新たなホルマント軌跡を求め、その値から合成に用いる予測係数を算出する。なお、ここで、第4以上のホルマント及び無声音部、有声音部は変更しない。

(d) 新たな予測係数と最初に求めた残差から合成音を発生する。

この第2の従来例においては、ホルマント周波数のダイナミクスを強調又は抑制するのみなので、声の明瞭性を改善することはできるが、第1の従来例と同様に、異なる個性を有する音声を発生することはできず、声質制御の自由度が小さいという問題点があった。

【0004】 さらに、ある特定話者からターゲット話者への声質の変換方法（以下、第3の従来例という。）が、M. Abe et al., “Voice Conversion through vector quantization”, Proc. ICASSP'88, pp. 655-658, 1988年に開示されている。この第3の従来例においては、いわゆるコードベクトル・マッピング手法に基づいて特徴パラメータのベクトル量を制御し、これらのマッピングは音声スペクトルに対する適切な拘束なしに学習データから計算されていたために、ターゲット話者への適切なマッピング関数を求めるときに、ターゲット話者による大量の発声データを必要とし、極めて大きな記憶装置を設ける必要があるという問題点があった。

【0005】 本発明の目的は以上の問題点を解決し、従来例に比較して非常に簡単な構成を有し、しかもより高い自由度を有して声質制御して音声合成することができる音声合成装置を提供することにある。

【0006】

【課題を解決するための手段】 本発明に係る請求項1記載の声質制御型音声合成装置は、複数の話者の音声スペクトルのデータを予め記憶する記憶手段と、上記記憶手段から複数の話者の音声スペクトルのデータを読み出し、所定の補間比率を用いて上記複数の話者の音声スペクトルのデータに対して内挿処理を実行して所定の声質

3

を有する音声スペクトルのデータを出力する処理手段と、上記処理手段から出力される音声スペクトルのデータに基づいて音声を合成して出力する音声合成手段とを備えたことを特徴とする。

【0007】また、請求項2記載の声質制御型音声合成装置は、請求項1記載の声質制御型音声合成装置において、さらに、ターゲット話者の音声に基づいて音声スペクトルのデータを抽出する特徴抽出手段と、上記特徴抽出手段によって抽出された音声スペクトルのデータが上記処理手段から出力される音声スペクトルのデータに近似するように上記補間比率を演算して上記処理手段に設定する演算手段とを備えたことを特徴とする。

【0008】さらに、請求項3記載の声質制御型音声合成装置は、請求項1記載の声質制御型音声合成装置において、さらに、上記補間比率を入力して上記処理手段に設定する入力手段を備えたことを特徴とする。

【0009】またさらに、請求項4記載の声質制御型音声合成装置は、請求項1、2又は3記載の声質制御型音声合成装置において、上記音声スペクトルのデータは、ケプストラム係数又は対数面積比であることを特徴とする。

【0010】

【作用】音声スペクトルを変更することにより、声質を制御するためには、音声のスペクトル構造及びそのダイナミクスに存在するある種の特徴を適切にモデル化し、そのモデル化に基づいたスペクトルの制御を行うことが望ましいと考えられる。しかしながら、音声スペクトルのモデル化を直接に、声門運動や声道形状に基づいて行った音声合成方式では、ホルマント合成で人手による精密な制御を行った数例を除き、自動的な手段を用いて高品質な合成音を発声させることはできていない。そこで、本発明者は、このような物理的モデルを直接的に用いる代わりに複数の人数の音声スペクトルそのものをノン・パラメトリックな音声スペクトルモデルとみなし、これをスペクトル制御の拘束条件として用いることを考えた。新しい音声スペクトルは、以下に詳細後述するように、複数の人数の音声スペクトルを線形に内挿することによって求める。

【0011】請求項1記載の声質制御型音声合成装置においては、上記処理手段は、上記記憶手段から複数の話者の音声スペクトルのデータを読み出し、所定の補間比率を用いて上記複数の話者の音声スペクトルのデータに対して内挿処理を実行して所定の声質を有する音声スペクトルのデータを出力し、次いで、上記音声合成手段は、上記処理手段から出力される音声スペクトルのデータに基づいて音声を合成して出力する。

【0012】また、請求項2記載の声質制御型音声合成装置においては、請求項1記載の声質制御型音声合成装置において、さらに、上記特徴抽出手段は、ターゲット話者の音声に基づいて音声スペクトルのデータを抽出

4

し、上記演算手段は、上記特徴抽出手段によって抽出された音声スペクトルのデータが上記処理手段から出力される音声スペクトルのデータに近似するように上記補間比率を演算して上記処理手段に設定する。これによって、上記ターゲット話者の音声に近似した音声を上記音声合成手段によって合成することができる。

【0013】さらに、請求項3記載の声質制御型音声合成装置においては、請求項1記載の声質制御型音声合成装置において、さらに、上記入力手段は、上記補間比率を入力して上記処理手段に設定する。従って、上記補間比率を変更して種々の声質を有する音声を上記音声合成手段によって発声させることができる。

【0014】またさらに、請求項4記載の声質制御型音声合成装置においては、請求項1、2又は3記載の声質制御型音声合成装置において、上記音声スペクトルのデータは、好ましくは、ケプストラム係数又は対数面積比である。

【0015】

【実施例】以下、図面を参照して本発明に係る実施例について説明する。図1は本発明に係る一実施例である声質制御型音声合成装置のブロック図である。

【0016】本実施例の音声合成装置は、複数の人数の音声スペクトルを線形に内挿する処理のために、予め音声情報メモリ10-1乃至10-Mにそれぞれ複数M人の話者の音声スペクトルのデータを格納し、ターゲット話者の音声を入力してその音声に近似するように、上記複数M人の話者の音声スペクトルを内挿混合することにより、上記ターゲット話者の音声の声質に近い音声を発声させる一方、上記内挿比率を変更することによって所定の異なる声質を有する音声を発声させることを特徴としている。ここで、内挿するためのスペクトル・パラメータとして、比較的良好な補間特性を有するケプストラム係数又は対数面積比を用い、パラメータの内挿比率を変更することにより、異なる声質を有する音声を合成して出力させる。

【0017】ターゲット話者の発声音声はマイクロホン1に入力されて音声信号に変換された後、特徴抽出部2に入力される。一方、音声情報メモリ10-1乃至10-Mにそれぞれ、複数M人の話者の音声スペクトルのデータを格納されている。ここで、音声スペクトルのデータは、音声スペクトルの振幅の時系列データ及び例えば16次のケプストラム係数の時系列データを含む。

【0018】特徴抽出部2は、入力された音声信号をA/D変換した後、例えばLPC分析を実行し、対数パワー、16次ケプストラム係数、 Δ 対数パワー及び16次 Δ ケプストラム係数を含む34次元の特徴パラメータを抽出する。抽出された特徴パラメータの時系列はバッファメモリ3を介して時間正規化部4に入力される。時間正規化部4は、ターゲット話者の発声音声のスペクトルと、上記音声情報メモリ10-1乃至10-Mに予め記

5

憶された複数M人の話者のスペクトルとの時間整合を、距離尺度としてケプストラム距離を用いてDTW (Dynamic time warping) 法により実行する。すなわち、ターゲット話者の発声音声の例えば単語又は文の時間長さは人及び時々により変化するので、当該ターゲット話者の発声音声のスペクトルのデータを、その単語又は文と同一の単語又は文に関する複数M人の話者の音声スペクトルの時間長さと同一となるように時間整合処理(時間正規化処理)を実行し、処理後のターゲット話者の音声スペクトルのデータは補間比率決定部5に出力される。

【0019】音声情報メモリ10-1乃至10-Mから読み出される複数M人の音声スペクトルのケプストラム係数データは補間比率決定部5に出力されるとともに、乗算器11-1乃至11-Mに出力される。乗算器11-1乃至11-Mはそれぞれ、入力された各人の音声スペクトルのケプストラム係数データと、補間比率決定部5から出力される補間比率 a_1, a_2, \dots, a_M とを乗算して加算器12に出力し、加算器12は入力されるデータを加算して、加算結果のデータをパラメータ時系列生成部6に出力する。すなわち、複数M個の乗算器11-1乃至11-Mと加算器12とによって音声スペクトルの内挿処理が実行される。

【0020】補間比率決定部5は、ターゲット話者への声質適応を行う場合、すなわち音声合成後のターゲット話者の発声音声に近似させる場合、ターゲット話者のスペクトルとスペクトルの内挿により生成したスペクトルとの間の距離が最小になるように内挿比率を決定する。具体的には、最適な補間比率 a_1, a_2, \dots, a_M を次の数1で示す関数Fの関数値を最小2乗法により演算して決定する。すなわち、ターゲット話者の発声音声のケプストラム係数値と、それに時間的に対応する予め格納された複数M人の音声スペクトルのケプストラム係数値との差の二乗が最小になるように、補間比率 a_1, a_2, \dots, a_M を求める。

【0021】

【数1】

$$F(a_1, a_2, \dots, a_M) = \sum_{i,j} (Y_{i,j} - y_{i,j})^2$$

ここで、

【数2】

$$y_{i,j} = a_1 \cdot x_{1,i,j} + a_2 \cdot x_{2,i,j} + \dots + a_M \cdot x_{M,i,j}$$

ただし、

【数3】

$$\sum_k a_k = 1$$

である。

【0022】ここで、 $Y_{i,j}$ と $y_{i,j}$ はそれぞれ、ターゲット話者と内挿により得られるスペクトルの第1フレームのj次ケプストラム係数を表わす。 $x_{k,i,j}$ は予め音声情報メモリ10-1乃至10-Mに格納されたk番目の話

6

者の第1フレームのj次ケプストラム係数を表わす。なお、音声合成後のターゲット話者の発声音声に近似させず、ターゲット話者とは異なる声質の音声合成する場合は、補間比率を適宜変更する。この場合、補間比率は、操作者がキーボード(図示せず。)を用いて補間比率決定部5に入力するように構成してもよい。

【0023】上記補間比率決定部5の処理の後に、パラメータ時系列生成部6は、加算器12から逐次出力される16次のケプストラム係数の時系列を取りまとめて、内蔵のバッファメモリに格納した後、そのデータを係数変換部7に出力する。係数変換部7は、入力された16次のケプストラム係数の時系列データに基づいて、そのデータを、公知の方法により、音声合成のためのピッチ、有聲/無聲切り換え、振幅及びフィルタ係数のデータに変換して、それぞれパルス発生器21とスイッチSWと振幅変更型増幅器23とフィルタ24とに出力する。

【0024】音声合成部20は、パルス発生器21と雑音発生器22とスイッチSWと振幅変更型増幅器23とフィルタ24とから構成される。パルス発生器21は、有声音の励振音源であって各ピッチ周期の開始時点で単位大きさのインパルスを発生して、スイッチSWを介して振幅変更型増幅器23に出力する。一方、雑音発生器22は、無声音の励振音源であって、無相関でかつ一様分布を有する標準偏差1と平均値0のランダム雑音を発生して、スイッチSWを介して振幅変更型増幅器23に出力する。従って、スイッチSWは有声音を発生するときパルス発生器21側に切り換える一方、無声音を発生するときは雑音発生器22側に切り換えられる。さらに、振幅変更型増幅器23は、入力される振幅情報に基づいて入力される信号の振幅を変更しかつ増幅してフィルタ24に出力する。そして、フィルタ24は、その伝達関数に対応するフィルタ係数を入力されるフィルタ係数に設定し、入力された信号を当該設定されたフィルタ係数でろ波した後、スピーカ30を介して出力する。

【0025】以上実施例において、音声合成後のターゲット話者の発声音声に近似させる場合、このスピーカ30からは、ターゲット話者の発声音声に近似した音声信号が出力される一方、音声合成後のターゲット話者の発声音声に近似させず、ターゲット話者とは異なる声質の音声合成する場合は、設定された補間比率に対応して声質制御された音声信号を出力させることができる。なお、後者の場合においては、ターゲット話者の発声音声に基づかず、補間比率決定部5以降の回路のみで構成してもよい。また、音声情報メモリ10-1乃至10-Mに格納される音声スペクトルのデータは予め時間正規化処理を実行されていることが好ましい。

【0026】以上の実施例において、スペクトル・パラメータとして、ケプストラム係数を用いているが、本発明はこれに限らず、PARCORパラメータ k_i から誘

導できる等価パラメータ集合を示す次の数4の対数面積比パラメータ g_i を用いてもよい。この場合、補間比率決定部5は、音声合成後のターゲット話者の発声音声に近似させるとき、好ましくは、数1を用いて演算した補*

$$g_i = \log \left[(1 - k_i) / (1 + k_i) \right], \quad 1 \leq i \leq p$$

【0027】本発明者は、本実施例の装置を用いてシミュレーションを実行して、2話者の中で補間が良好に行えるかどうかを調べた。音質に関しては、ケプストラム係数と対数面積比のどちらのパラメータを用いた場合も比較的良好に補間して内挿することができ、補間比率を変化することにより声質が一方の話者から他方の話者に安定に徐々に変化していくことを確認した。

【0028】

【発明の効果】以上詳述したように本発明によれば、複数の話者の音声スペクトルのデータを予め記憶する記憶手段と、上記記憶手段から複数の話者の音声スペクトルのデータを読み出し、所定の補間比率を用いて上記複数の話者の音声スペクトルのデータに対して内挿処理を実行して所定の声質を有する音声スペクトルのデータを出力する処理手段と、上記処理手段から出力される音声スペクトルのデータに基づいて音声を合成して出力する音声合成手段とを備えたので、従来例に比較してより小さい記憶容量を有する記憶装置を用い、より簡単な構成の回路を用いて、より自由度が高い声質の制御が可能になり、より多様な声質を有する音声を合成することができ

る。
【0029】さらに、ターゲット話者の音声に基づいて音声スペクトルのデータを抽出する特徴抽出手段と、上記特徴抽出手段によって抽出された音声スペクトルのデータが上記処理手段から出力される音声スペクトルのデータに近似するように上記補間比率を演算して上記処理

*間比率を初期値として、さらに、非線形降下法を用いてケプストラム距離の低減を行う。

【数4】

手段に設定する演算手段とを備えたので、ターゲット話者が発声した少量の音声を入力として、声質変換のターゲットとする話者の声質を近似する音声を発生することができるという特有の利点がある。

【図面の簡単な説明】

【図1】 本発明に係る一実施例である声質制御型音声合成装置のブロック図である。

【符号の説明】

- 1…マイクロホン、
- 2…特徴抽出部、
- 3…バッファメモリ、
- 4…時間正規化部、
- 5…補間比率決定部、
- 6…パラメータ時系列生成部、
- 7…係数変換部、
- 10-1乃至10-M…音声情報メモリ、
- 11-1乃至11-M…乗算器、
- 12…加算器、
- 20…音声合成部、
- 21…パルス発生器、
- 22…雑音発生器、
- 23…利得変更型増幅器、
- 24…フィルタ、
- 30…スピーカ、
- SW…スイッチ。

152

